

libraries were spotted on nylon membrane filters and screened with oligonucleotide probes (e.g., 7-mers) to obtain signature sequences. The inserts for the cDNA libraries from which the sequences were obtained were amplified with PCR using primers specific for the vector sequences which flank the inserts, or isolated from plasmid preparations. The 5' sequence of the amplified inserts was then deduced using the reverse M13 sequencing primer in a typical Sanger sequencing protocol, as well as internal primers in both the forward and reverse direction. In some cases RACE (Random Amplification of cDNA Ends) was performed to further extend the sequence in the 5' direction. In all cases all of a signature cluster was sequenced to generate overlapping clones to assemble the contigs. Chromatograms were base called and assembled using a software suite from University of Washington, Seattle containing three applications designated PHRED, PHRAP, and CONSED. The sequences for the resulting contigs for the 782 gene family are designated as 782 SEQ ID NO: 1-10,451 and are provided in the attached Sequence Listing. inserts was then deduced in a typical Sanger sequencing protocol. The inserts of the library were, amplified with PCR using 5 primers specific for vector sequences which flank the inserts.

The contigs were assembled using an EST sequence as a seed. Then a recursive algorithm was used to extend the seed EST into an extended assemblage, by pulling additional sequences from different databases (*i.e.*, Hyseq's database containing EST sequences, dbEST version 114, gb pri 114, and UniGene version 101) that belong to this assemblage. The algorithm terminated when there was no additional sequences from the above databases that would extend the assemblage. Inclusion of component sequences into the assemblage was based on a BLASTN hit to the extending assemblage with BLAST score greater than 300 and percent identity greater than 95%.

The nearest neighbor result for the assembled contig was obtained by a FASTA version 3 search against Genpept release 114, using FASTXY algorithm. FASTXY is an improved version of FASTA alignment which allows in-codon frame shifts. The nearest neighbor result showed the closest homologue for each assemblage from Genpept (and contains the translated amino acid sequences for which the assemblage encodes). The nearest neighbor results for 782 SEQ ID NO: 1-10,451 are shown in Table 2, and identified as Table2(782).doc on the enclosed compact disc.

## 6.6 The 784 Gene Family

### Novel Contigs

Table 3 (identified as Table3(784).doc on the enclosed CD) sets forth the novel predicted polypeptides (including proteins) encoded by the novel polynucleotides (784 SEQ ID NO: 1-10,289) of the present invention, and their corresponding nucleotide locations to each of 748 SEQ ID NO: 1-10,289. Table 3 also indicates the method by which the polypeptide was predicted. Method A refers to a polypeptide obtained by using a software program called FASTY (available from <http://fasta.bioch.virginia.edu>) which selects a polypeptide based on a comparison of translated novel polynucleotide to known polypeptides (W.R. Pearson, Methods in Enzymology, 183: 63-98 (1990), incorporated herein by reference). Method B refers to a polypeptide obtained by using a software program called GenScan for human/vertebrate sequences (available from Stanford University, Office of Technology Licensing) that predicts the polypeptide based on a probabilistic model of gene structure/compositional properties (C. Burge and S. Karlin, J. Mol. Biol., 268: 78-94 (1997), incorporated herein by reference). Method C refers to a polypeptide obtained by using a Hyseq proprietary software program that translates the novel polynucleotide and its complementary strand into six possible amino acid sequences (forward and reverse frames) and chooses the polypeptide with the longest open reading frame. When the predicted beginning nucleotide of Table 3 is a higher number than the predicted end nucleotide of Table 3, then the amino acid sequence is derived from the complementary strand of the indicated SEQ ID NO. The locations of the predicted beginning and end nucleotides correlate to the nucleotide sequence of the indicated SEQ ID NO., not its complementary strand.

The isolated polypeptides of the invention include, but are not limited to, a polypeptide comprising any of the amino acid sequences set forth in Table 3 or from six frame translations of 784 SEQ ID NO: 1-10,289; or the corresponding full length or mature protein. One of skill in the art could determine the corresponding amino acid sequence using techniques well known in the art to translate and analyze all possible six frames. Polypeptides of the invention also include polypeptides with biological activity that are encoded by (a) any of the polynucleotides having a nucleotide sequence set forth in the 784 SEQ ID NO: 1-10,289; or (b) polynucleotides that hybridize to the

Table 4 (identified as Table4(784).doc on the enclosed CD) shows the various tissue sources of the EST sequences from Hyseq's database which were used to assemble the contigs or nucleic acids of the present invention (identified by 784 SEQ ID NO: 1-10,289).

5       The nearest neighbor result for the assembled contig was obtained by a FASTA version 3 search against Genpept release 114, using FASTXY algorithm. FASTXY is an improved version of FASTA alignment which allows in-codon frame shifts. The nearest neighbor result showed the closest homologue for each assemblage from Genpept (and contains the translated amino acid sequences for which the assemblage encodes). The  
10       nearest neighbor results for 784 SEQ ID NO: 1-10,289 are shown in the Table 5, infra.

## 6.7    **The 785 Gene Family**

### Novel Nucleic Acid Sequences Obtained From Various Libraries

A plurality of novel nucleic acids were obtained from cDNA libraries prepared  
15       from various human tissues and in some cases isolated from a genomic library derived from human chromosome using standard PCR, SBH sequence signature analysis and Sanger sequencing techniques. The inserts of the library were amplified with PCR using primers specific for the vector sequences which flank the inserts. Clones from cDNA  
libraries were spotted on nylon membrane filters and screened with oligonucleotide  
20       probes (*e.g.*, 7-mers) to obtain signature sequences. The clones were clustered into groups of similar or identical sequences. Representative clones were selected for sequencing.

In some cases, the 5' sequence of the amplified inserts was then deduced using a typical Sanger sequencing protocol. PCR products were purified and subjected to  
fluorescent dye terminator cycle sequencing. Single pass gel sequencing was done using  
25       a 377 Applied Biosystems (ABI) sequencer to obtain the novel nucleic acid sequences. In some cases RACE (Random Amplification of cDNA Ends) was performed to further extend the sequence in the 5' direction.

The novel contigs of the invention were assembled from sequences that were obtained from a cDNA library by methods described above, and in some cases sequences  
30       obtained from one or more public databases. Chromatograms were base called and assembled using a software suite from University of Washington, Seattle containing three